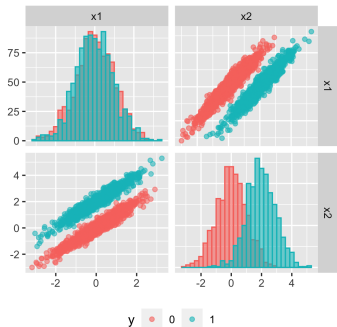


Introduction to Machine Learning

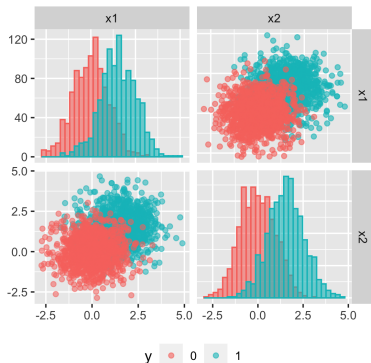
Feature Selection: Filter Methods (Examples and Caveats)



Learning goals

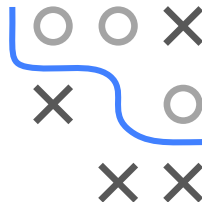
- Understand how filter methods can be misleading.
- Understand how filters can be applied and tuned.

FILTER METHODS CAN BE MISLEADING



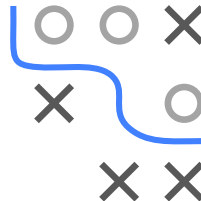
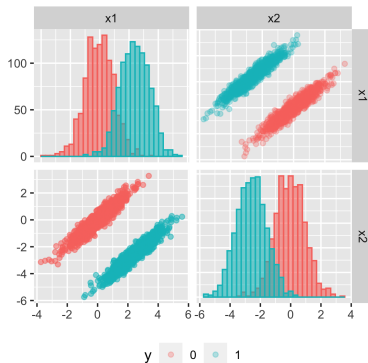
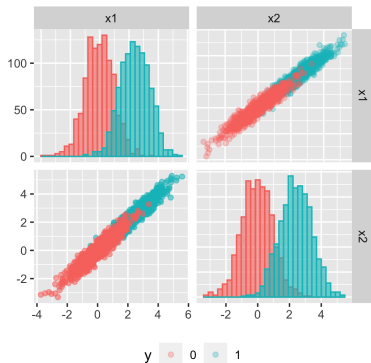
ρ_{ACC} of log. reg. classifier with:

- feature x_1 : 0.76
- feature x_2 : 0.78
- both features: 0.85



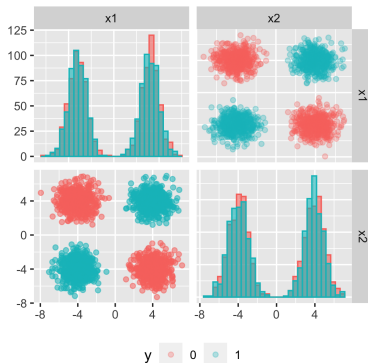
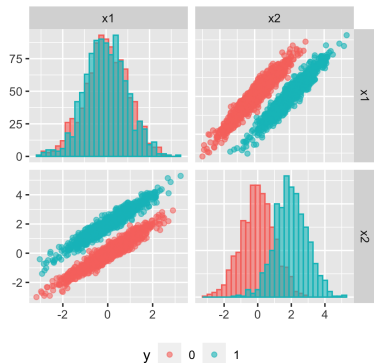
IG from presumably redundant variables. 2 class problem with i.i.d. variables. Each class has Gaussian distribution with no covariance. While filter methods suggest redundancy, combination of both vars yields improvement, showing i.i.d. vars are not truly redundant. For further details, see [► Guyon and Elisseeff, 2003](#).

FILTER METHODS CAN BE MISLEADING



Intra-class covariance. In projection onto the axes, distribution of two variables are same as before. Left: Class conditional distribution have high cov. in direction of the line of the two class centers. Right: Class conditional distr. have high cov. in direction perpendicular to line of two class centers. Important separation gain is obtained by using both variables.

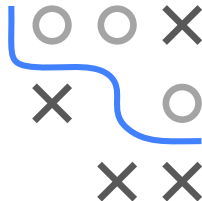
FILTER METHODS CAN BE MISLEADING



Variable useless by itself can be useful together with others. Left: One var has completely overlapping class conditional densities. Still, jointly with other variable separability can be improved. Right: XOR-like chessboard problem. Classes consist of “clumps” s.t. projection on the axes yields overlapping densities. Single vars have no separation power, only used together.

- 1 Calculate filter score for each feature x_j
- 2 Rank features according to score values.
- 3 Choose \tilde{p} best features
- 4 Train model on \tilde{p} best features.

- Could be prescribed by the application
- Eyeball estimation: read from filter plots
- Treat as hyperparameter and tune in a pipeline, based on resampling



USING FILTER METHODS

Advantages:

- Easy to calculate
- Typically scales well with the number of features p
- Generally interpretable
- Model-agnostic

Disadvantages:

- Univariate analyses may ignore multivariate dependencies
- Redundant features will have similar weights
- Ignores the learning algorithm

